

# Automatic Vocal Mixing Using Convolutional Neural Networks

kelian Li  
Center for Music Technology  
Georgia Institute of Technology  
Atlanta, USA  
kli421@gatech.edu

Alexander Lerch  
Center for Music Technology  
Georgia Institute of Technology  
Atlanta, USA  
alexander.lerch@gatech.edu

*Abstract*—abstract

abstract

abstract

abstract

abstract

abstract

abstract

abstract

abstract

abstract

*Index Terms*—automatic mixing, intelligent music production

## I. INTRODUCTION

Music mixing is the process of blending multi-track recordings through manipulating audio effects, and the mixing quality is ultimately judged by the preference of listeners. Automatic mixing aims to reach a similar level of quality by computational systems without human interventions [1]. Automatic mixing research can not only increase the efficiency, but also help understand mixing from both scientific and artistic perspectives which may lead to innovations in music production.

This work intends to tackle the problem of vocal and backing track mixing by deep learning. Previous research attempted to create general-purposed multi-track mixing systems which has not yet succeeded. This study focuses on vocal and backing track mixing which has not been addressed in previous research but with wide applications in professional music production and public music participation.

Derived from hip-hop culture, a considerable amount of music nowadays is created with pre-made accompaniment tracks [2]. Unlike the traditional song-writing and music production where these two processes go hand in hand, the communication between producers and artists is no longer a necessary process. This new norm of music making expanded out of the hip-hop community and has been adapted by the broader music industry as shown in the success of beat-selling platforms like BeatStars [3]. Automatic vocal mixing can accelerate the mixing process and reduce the listening fatigue of the mixing engineers.

Another potential application is karaoke. Mobile karaoke apps enable users to sing and share their performances on social media [4]. Most karaoke mobile apps have some basic

vocal mixing features, but general users are unable to produce very good mixes without much knowledge of music mixing.

## II. RELATED WORK

Early automatic mixing methods are mostly rule-based expert systems. These systems incorporate some common practices by human mixing engineers into the algorithms that can perform similar procedures. Reiss [5] firstly introduced the concept of adaptive audio effects, and the key idea is that the audio effect parameters are modified by analysis of input audio through feature extraction and mapping. More comprehensive overview of rule-based automatic mixing systems can be found in the work by De Man et al. [6], Moffat and Sandler [1], and De Man et al. [7].

As described in the book chapter by Moffat [8], the major limitation of the data-driven automatic mixing is the data collection of mixing parameters. One solution is to let users create training data [9] [10]. Another common solution is to adapt the reverse engineering approach [11] [12]. Martínez Ramírez et al. [13] propose a method to extract the mixing parameters for arbitrary black-box audio effects by gradient approximation. All of these approaches require the raw multi-track recordings and the corresponding professional mix-downs.

End-to-end audio transformation provides a path without extracting the mixing parameters. Ramírez et al. [14] demonstrates an end-to-end mixing system for drum mixing with Wave-U-Net [15]. Based on the previous research of black-box audio effects [16], Ramírez et al extend the end-to-end music mixing to wet processed multitrack training data. Both systems can achieve mixing results that are indistinguishable from the professional mix according to their listening tests [17]. The most latest work shows that the same method can perform end-to-end music mixing style transfer [18].

One obvious drawback of the end-to-end framework is that it prevents further parameter control by humans. Steinmetz et al. [19] incorporate differentiable digital signal processing [20] to extract the mixing parameter data, and thus their system can produce human-readable mixing parameters. However, the implementation is very challenging, and the listening test result indicates that the system performance cannot match that of audio engineers.

Outside of academia, companies have begun adapting automatic mixing into their products. The list includes AYAIC<sup>1</sup>, RoEx<sup>2</sup>, Neutron<sup>3</sup>, Balancer<sup>4</sup>, and Faders<sup>5</sup>. In the task of automatic mastering, Landr<sup>6</sup>, eMastered<sup>7</sup>, Ozone<sup>8</sup>, BandLab<sup>9</sup>, and SoundCloud<sup>10</sup> are considered successful for achieving the commercially usable quality.

### III. METHODOLOGY

The audio effects applied to the lead vocals are in the sequence of equalization, compression, reverberation and level balance. In order to enable the mixing engineer to have the same signal flow, we choose the open-sourced audio plugins MultiEQ, OmniCompressor and FdnRverb by IEM<sup>11</sup>. Pedalboard<sup>12</sup> is used to load the plugins in the Python environment.

A baseline system and a deep learning system are implemented in this work. The baseline system combines both data analysis and rule-based approaches from previous studies. The deep learning regression model predicts either the direct or intermediate mixing parameters from the input signal.

#### A. Data Collection and Post-processing

Both the data analysis method and the deep learning method rely on the dataset. Using music source separation [21], the vocals and the backing tracks are extracted from 21936 clips of 30-second audio in a subset of the Million Song Dataset [22]. This data is used in level balance and compression. For equalization and reverberation, the dataset is the 100 full-length songs from the train set of MUSDB18HQ [23].

**Level Balance**—Relative integrated loudness [24] between the vocal and the backing track is the chosen loudness measure which can be easily transferred in the gain changes. The separated Million Song Dataset [22] provides the necessary data for the baseline and the deep learning model. To reduce the influence of the outliers, the songs with vocal relative loudness outside of [-10dB, 6dB] are removed from the dataset.

**Compression**—EBU loudness range [25] of the vocal signal is the dynamic measure in this work. The data is also collected from the separated Million Song Dataset [22]. The songs with vocal loudness range outside of [5dB, 30dB] are removed. An iterative process optimizes the vocal loudness range to 16.36dB  $\pm$  1dB by modifying the compression

threshold and compression ratio. The compression threshold and the compression ratio are initialized to 0dB and 3:1. For each iteration, the compression threshold is reduced by 3dB, and the compression ratio is increased by 0.1. The process ends if the loudness range is below 16.36dB, or the current compression threshold is lower than -35dB. The attack time and the release time are set to 30ms and 150ms. This process follows the general concept of adaptive audio effects proposed by Reiss [5].

**Equalization**—The purely rule-based baseline equalization described below does not require any data. The deep learning equalization model uses the 100 songs from the train set of MUSDB18HQ [23], but the data collection method is different from level balance and compression. The main idea is to create "raw" vocals by applying randomized EQ on the processed vocals. 4 out of 9 center frequencies in [63, 125, 250, 500, 1K, 2K, 4K, 8K, 16K] Hz are randomly chosen. Uniform distributed gain changes between [-15, 15] dB are then applied on the processed vocals using the MultiEQ plugin with the Q factor of 1. 73900 "raw" vocals in total are created by the 100 processed vocal tracks.

**Reverberation**—Blind reverb estimation is an unsolved problem in public research [26]. In this work, a commercial plugin Chameleon<sup>13</sup> is utilized to directly extract the estimated reverb impulse responses from the wet vocal signal in the 100 songs from the train set of MUSDB18HQ [23]. The next step is to match the reverb parameters with the impulse responses. As suggested in [27] and [28], genetic algorithms can perform impulse response matching efficiently. The genetic algorithm Implementation from PyGAD [29] helps match the FdnRverb plugin parameters with each estimated impulse response. The fitness function minimizes the magnitude spectrogram difference between the targeted and the generated impulse responses. In the end, the reverb parameters are approximated from the wet vocal tracks.

#### B. Baseline System

**Level Balance**—In an analysis of 64 mixes, De Man et al. [30] find that the relative loudness of the vocal track and the full mix loudness is about -2.7dB. From the 21936 clips in the separated Million Song Dataset [22], we discover that the integrated loudness of the vocals is 1.77dB lower than the backing tracks on average. The baseline system sets the relative integrated loudness of -1.77dB between the vocals and the backing tracks.

**Compression**—The average loudness range of the processed vocals in the separated Million Song Dataset [22] is 16.36dB. which is the desired vocal loudness range in the baseline system. The iterative post-processing described above sets the compression parameters which allows the final loudness range to be closed to 16.36dB.

<sup>1</sup>AYAIC, <https://www.ayaicinc.com/>

<sup>2</sup>RoEx, <https://www.roexaudio.com>

<sup>3</sup>Neutron, <https://www.izotope.com/en/products/neutron/features/mix-assistant.html>

<sup>4</sup>Balancer, <https://www.sonible.com/balancer/>

<sup>5</sup>Faders, <https://faders.io/>

<sup>6</sup>Landr, <https://www.landr.com/>

<sup>7</sup>eMastered, <https://emastered.com/>

<sup>8</sup>Ozone, <https://www.izotope.com/en/products/ozone/features/master-assistant.html>

<sup>9</sup>BandLab, <https://www.bandlab.com/mastering>

<sup>10</sup>SoundCloud, <https://community.soundcloud.com/mastering-on-soundcloud>

<sup>11</sup>IEM, <https://plugins.iem.at/docs/pluginDescriptions/>

<sup>12</sup>Pedalboard, <https://github.com/spotify/pedalboard>

<sup>13</sup>Chameleon, <https://www.accentize.com/chameleon/>

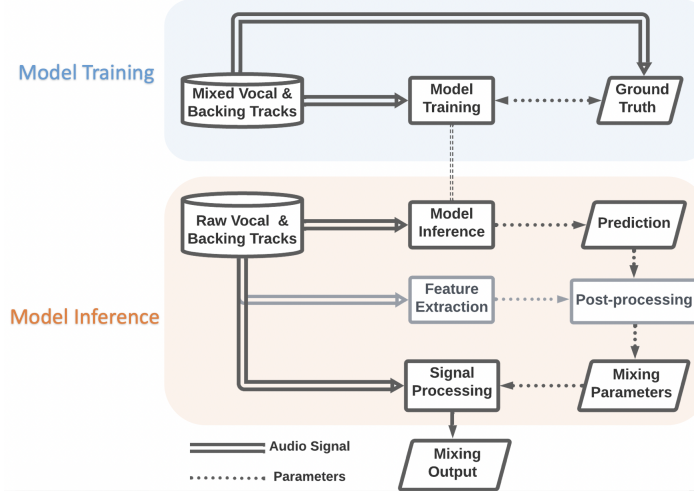


Fig. 1: Overview of the deep learning regression approach

**Equalization**—The baseline equalization is achieved by the rule-based frequency masking reduction method proposed by Hafezi and Reiss [31]. Masking is defined as the process by which the audibility of one sound (the maskee) is reduced by the presence of another sound (the masker). Equalizers are commonly used to reduce the spectral masking in mixing. Our implementation follows [31] with some modifications.

In short, the algorithm firstly identifies the essential and nonessential band of the two signals, and then applies peak filters when one essential band is masked when the same frequency band is nonessential in the maskee signal. The vocal and the backing track can both be the masker and the maskee, but the filter is only applied to the vocal. The incoming signal is firstly gated to reduce noise. The center frequencies of 9 octave bands in use are [63, 125, 250, 500, 1K, 2K, 4K, 8K, 16K] Hz, and the lowest 32.5Hz in the original implementation is eliminated. The RMS differences of the 9 bands are normalized by subtracting their mean. The top 6 bands of the vocal and the top 2 bands of the backing track are considered essential. Up to 4 frequency bands are selected where the filters are applied. The Q factor is changed from 2 to 1.

**Reverberation**—The automatic reverberation combines the data analysis and rule-based approaches. As suggested by Moffat et al., the reverb time may be correlated with the tempo [32]. In our modified equation where the tempo is linearly mapped to the reverb time,

$$RT = -\frac{1}{45}T + \frac{40}{9} \quad (1)$$

RT is the reverb time, and T is the tempo.

Other reverb parameters set to mean values extracted from the train set of MUSDB18HQ [23].

### C. deep learning System

Four independent convolutional neural network regression models are trained to predict the direct or intermediate mixing

parameters for each audio effects from the mel-spectrogram of the vocal and the backing track. The process of converting intermediate parameters to direct parameters is the same as it in the baseline system. The structures of the four models are the same except the last output layers.

To avoid the influence between audio effects, the signal is processed by the audio effect before being fed into the next model. The sequence of the four models is the same as the signal chain: equalization, compression, reverberation, and level balance. These models are trained to minimize the mean-squared-error between the predictions and the ground truth. During the inference, the output predictions are averaged over the mel-spectrogram frames of the entire song.

The vocal relative loudness and vocal loudness range are the ground truth of the level balance model and the compression model. The ground truth for the reverb model are the 10 parameters of the reverb plugin. Thus the model is a multi-output regression model. Some of these reverb parameters are empirically more important than the others. To reflect this difference, the loss function of the reverb model is

$$L_{reverb} = 10 \cdot l_{DW} + 5 \cdot l_{RT} + 5 \cdot l_{RS} + l_{rest} \quad (2)$$

where  $DW$  is dry-wet ratio,  $RT$  is reverb time,  $RS$  is room size, and  $l_{rest}$  is the mean-squared-error of the rest of the parameters.

Due to the enormous amount of the training data generated from applying random EQ settings to the vocal tracks, the mel-spectrograms of each altered song is concatenated to only one mel-spectrogram. Thus, the model output does not need to be averaged because there will be only one set of output.

The ground truth of the EQ model are the gain values in the 9 frequency bands. As mentioned in the data collection section, the peak filters are applied on 4 out of 9 center frequencies. This leads to an unbalanced training data where 5 of 9 the gain values are 0dB because they are unchanged. To tackle that, the loss function has to be altered to preventing the model

	Level (dB)	loudness range (dB)	EQ gain (dB)	DW	RT	FT	RS	LCF	LCG	LCQ	HCF	HCG	HCQ
CNN	<b>1.64</b>	<b>2.63</b>	4.48	<b>6.13</b>	<b>1.006</b>	<b>0.401</b>	<b>7.30</b>	41.55	<b>3.62</b>	<b>0.355</b>	<b>3390</b>	<b>5.01</b>	<b>0.14</b>
mean	2.13	2.88	<b>3.33</b>	7.16	1.007	0.403	7.31	<b>40.26</b>	3.65	0.360	3462	5.35	0.17

from outputting zeros only. The 9 frequency bands are split into two groups depending on whether the ground truth gain values are changed or not. The modified loss function is

$$L_{EQ} = l_{changed} + 0.1 \cdot l_{unchanged} \quad (3)$$

where  $l_{changed}$  is the L2 loss of the 4 frequency bands where the ground truth values are changed, and  $l_{unchanged}$  is the loss of the 5 frequency bands where the ground truth values are zeros. During the inference, the top 4 frequency bands which have the largest absolute gain change prediction are selected.

1) *model input*: The input feature is the mel-spectrogram of the mono vocal and the backing track at the sample rate of 44100Hz. The FFT size is 2048 with the hop size of 1024. The frames with the vocal's integrated loudness below -30dB are removed. The number of mel filterbanks is 128. Each input feature block includes 64 mel-spectrum temporal steps, and each block represents around 1.51 seconds. The size of the input block is (2, 128, 64). Each input block channel is independently normalized to its maximum value. The neural network architecture is shown below.

Conv2d (input channels = 2, output channels = 8)
BatchNorm
ReLU
MaxPool2d
Conv2d (input channels = 8, output channels = 16)
BatchNorm
ReLU
MaxPool2d
Conv2d (input channels = 16, output channels = 32)
BatchNorm
ReLU
MaxPool2d
Conv2d(input channels = 32, output channels = 64)
BatchNorm
ReLU
MaxPool2d
Dropout (p = 0.3)
MLP (input channels = 768)

For all the convolution layers, the kernel size is 3, the stride is 1, and no zero-padding applied. The kernel size of the max pooling layers is 2.

#### IV. EXPERIMENT

##### A. Experimental Setups

The ground truth data is normalized between [0., 1.]. The Adam optimizer [33] is used for training. The batch size is 25 except for the EQ model. In level balance, the learning rate is 1e-3. In compression, the learning rate is 1e-4. In equalization, the learning rate is 1e-3 and the batch size is

30. In reverberation, the learning rate is 1e-5, and the weight decay is 1e-5.

The training set for level balance and compression is the 21936 clips of 30-second audio from the source-separated Million Song Dataset [22]. For the equalization and reverberation models, the training data is the 100 songs in training set of MUSDB18HQ [23]. The 48 songs in the test set of MUSDB18HQ is used for validation of all four audio effects.

##### B. Objective Evaluation

The validation result in mean-absolute-error is shown in the table above. The mean prediction is what the baseline system uses except for equalization and reverb time. The EQ gain error is the mean absolute error over the 9 frequency bands. DW, RT, FT, RS, LCF, LCG, LCQ, HCF, HCG, and HCQ are the reverb parameters in the FdnVerb plugin and each of them stands for dry-wet ratio (%), reverb time (s), Fade-in Time (s), Room Size, Low Cut Frequency (Hz), Low Cut Gain (dB), Low Cut Q factor, High Cut Frequency (Hz), High Cut Gain (dB), and High Cut Q factor respectively. Room Size ranges from 1 to 30.

In the tasks of level balance and compression, the CNN model can outperform the baseline system. On reverb parameter prediction, the CNN model can only perform slightly better than the baseline model on most of the parameters. The deep learning model successfully converges during the training, but the validation error shows that the model lacks the ability of generalization.

##### C. Listening Test

Due to the subjective nature of mixing, a listening test is conducted to evaluate the mixing quality of the two proposed systems. Besides the mixes by the baseline system and the deep learning system, the test participants are presented with a human mix as the high anchor and a random mix as the low anchor. The human mix uses the same audio plugins from the automatic mixing systems and keeps constant parameters such as the compression attack time and the peak filter's Q factors untouched. Uniform random mixing parameters are created for reverb parameters, compression threshold and ratio, and EQ gains. The vocal relative loudness is randomly generated from the range of [-6, 6] dB in order to make the vocals audible. The mixing limitations are kept the same way as they are in the human mixes. All mixes are loudness normalized to -28dBFS.

Apart from many previous mixing listening tests, we aim to include more songs in shorter length within 20-30 minutes of listening time. The participants are asked to listen to 10 seconds of 15 songs in 4 mix versions. The received feedback indicates that this is a workable method to gather more testing data if the mixes do not include automation and song structures. Amateur performances and low quality recordings are intentionally chosen for this listening test to reflect the

wider practice of music making. We also want to understand the effectiveness of audio mixing in non-professional settings. The 15 songs are selected from DAMP-VSEP [34] which is a real-world dataset karaoke performances on the Smule app. These songs are in different languages and different music styles. The participants rate the mixing quality in the likert scale of "very Poor", "Poor", "Fair", "Good", and "Very Good". Two examples of "Very Good" and "Very Poor" mixes are given before the official test begins. Beyond rating the overall mixing quality, the participants also rate the use of each audio effects in the same scale. Among the 30 online participants, 24 of them claim to have experiences in audio mixing. 10 participants identify themselves as professional audio engineers.

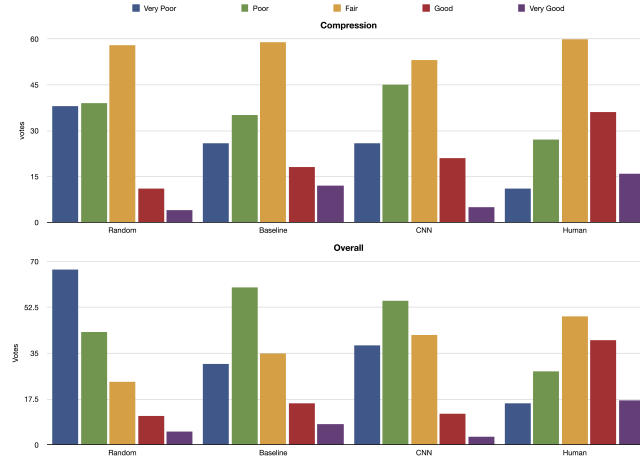


Fig. 2: The compression rating does not align with the overall rating

## V. DISCUSSION

(1) the lack of generalization ability of this model; (2) not enough training data; (3) improper input features; (4) improper ground truth representation; (5) the challenge of the task itself; (6) code errors.

(1): the lowest validation loss throughout the training is at the 2nd epoch hints that the model may have overfitted early. The solution can be increasing the complexity and the size of the model while employing more regularization techniques. (2): datasets of professional or semi-professional mixes like MedleyDB [35], iKala [36], ccMixer [37], and Rock band [38]. Source separation tools can also be used to extend the dataset from almost unlimited published recordings. (3): The current input feature represents the audio of 1.51 seconds. The segment may be too short, and features other than mel-spectrogram can also be taken into consideration. (4): The current ground truth is the relative loudness between the vocal and the backing track of the full song measured in dB. This representation can be changed into the local relative loudness of the audio segment. The loudness in dB can also be converted to amplitude or scaled down to [0, 1] by min-max normalization. (5): Neither the training set or the validation set represents the best mixing choice since *best* cannot even be defined in a creative task. The failure of the deep learning model is expected, and it is still meaningful result to understand automatic mixing. However, the model output is not yet processed and create a mix. The averaged prediction over the song and the listening test may show a different result. (6): The code errors in training the model and data processing must be eliminated before evaluating the model. A trial run of this model on a separated task may prove the reliability of both the code and the model.

### A. limitation

poor reverb algorithm

## VI. FUTURE WORK

### VII. CONCLUSION

A rule-based system for vocal and backing track automatic mixing is implemented. Based on literature and data analysis, the system can perform level balance, compression, equalization, and reverberation. A listening test will be conducted to evaluate the effectiveness of this system. A trial experiment of using convolution neural network to perform level balance shows that the model is not learning. Future work is needed to improving the data-driven level balance model.

### VIII. ACKNOWLEDGMENTS

We would like to appreciate Ben Holst for providing the human mixes in the listening test. We also thank all the participants in the listening test and their valuable feedback.

### REFERENCES

- [1] D. Moffat and M. B. Sandler, "Approaches in intelligent music production," in *Arts, Multidisciplinary Digital Publishing Institute*, vol. 8, 2019.
- [2] J. G. Schloss, *Making beats: The art of sample-based hip-hop*. Wesleyan University Press, 2014.
- [3] M. Stassen, *Having paid out \$150m to creators, BeatStars launches Sony Music Publishing-backed publishing service*, 2021. [Online]. Available: <https://www.musicbusinessworldwide.com/having-paid-out-150m-to-creators-beatstars-launches-sony-music-publishing-backed-publishing-service/>.
- [4] X. Zhou and F. Tarocco, *Karaoke: The global phenomenon*. Reaktion Books, 2013.
- [5] J. D. Reiss, "Intelligent systems for mixing multichannel audio," in *2011 17th International Conference on Digital Signal Processing (DSP)*, IEEE, 2011.
- [6] B. De Man, J. Reiss, and R. Stables, "Ten years of automatic mixing," in *3rd Workshop on Intelligent Music Production*, vol. 15, 2017.

- [7] B. De Man, R. Stables, and J. D. Reiss, *Intelligent Music Production*. Audio Engineering Society Presents. Taylor Francis, 2019.
- [8] D. Moffat, "Ai music mixing systems," in *Handbook of Artificial Intelligence for Music*, Springer, 2021, pp. 345–375.
- [9] E. T. Chourdakis and J. D. Reiss, "Automatic control of a digital reverberation effect using hybrid models," in *Dereverberation and Reverberation of Audio, Music, and Speech Conference*, Audio Engineering Society, 2016.
- [10] E. T. Chourdakis and J. D. Reiss, "A machine-learning approach to application of intelligent artificial reverberation," *Journal of the Audio Engineering Society*, vol. 65, no. 1/2, 2017.
- [11] D. Barchiesi and J. Reiss, "Reverse engineering of a mix," *Journal of the Audio Engineering Society*, vol. 58, no. 7/8, 2010.
- [12] J. T. Colonel and J. Reiss, "Reverse engineering of a recording mix with differentiable digital signal processing," *The Journal of the Acoustical Society of America*, vol. 150, no. 1, 2021.
- [13] M. A. Martínez Ramírez, O. Wang, P. Smaragdis, and N. J. Bryan, "Differentiable signal processing with black-box audio effects," in *ICASSP*, IEEE, 2021.
- [14] M. Martinez Ramirez, D. Stoller, and D. Moffat, "A deep learning approach to intelligent drum mixing with the wave-u-net," Audio Engineering Society, 2021.
- [15] D. Stoller, S. Ewert, and S. Dixon, "Wave-u-net: A multi-scale neural network for end-to-end audio source separation," in *Proceedings of ISMIR*, 2018.
- [16] M. A. Martinez Ramirez, E. Benetos, and J. D. Reiss, "Deep learning for black-box modeling of audio effects," *Applied Sciences*, vol. 10, no. 2, p. 638, 2020.
- [17] M. A. Martinez-Ramirez, W.-H. Liao, G. Fabbro, S. Uhlich, C. Nagashima, and Y. Mitsufuji, "Automatic music mixing with deep learning and out-of-domain data," in *Proceedings of ISMIR*, 2022.
- [18] J. Koo, M. A. Martinez-Ramirez, W.-H. Liao, S. Uhlich, K. Lee, and Y. Mitsufuji, "Music mixing style transfer: A contrastive learning approach to disentangle audio effects," *arXiv preprint arXiv:2211.02247*, 2022.
- [19] C. J. Steinmetz, J. Pons, S. Pascual, and J. Serrà, "Automatic multitrack mixing with a differentiable mixing console of neural audio effects," in *ICASSP*, IEEE, 2021.
- [20] J. Engel, C. Gu, A. Roberts, *et al.*, "DDSP: Differentiable digital signal processing," in *International Conference on Learning Representations*, 2019.
- [21] A. Défossez, "Hybrid spectrogram and waveform source separation," in *Proceedings of the ISMIR 2021 Workshop on Music Source Separation*, 2021.
- [22] T. Bertin-Mahieux, D. Ellis, B. Whitman, and P. Lamere, "The million song dataset," in *Proceedings of ISMIR*, 2011.
- [23] Z. Rafii, A. Liutkus, F.-R. Stöter, S. I. Mimilakis, and R. Bittner, *MUSDB18-HQ - an uncompressed version of musdb18*, Dec. 2019. [Online]. Available: <https://doi.org/10.5281/zenodo.3338373>.
- [24] EBU-Recommendation, "Loudness normalisation and permitted maximum level of audio signals," 2011.
- [25] E. T. 3342, "Loudness range: A measure to supplementebu r 128 loudness normalisation," 2016.
- [26] A. Sarroff and R. Michaels, "Blind arbitrary reverb matching," in *Proceedings of DAFx*, 2020.
- [27] A. Cunningham, K. McNally, and P. Driessen, "Automatic impulse response matching for reverb plugins," in *Audio Engineering Society Convention 150*, Audio Engineering Society, 2021.
- [28] S. Heise, M. Hlatky, and J. Loviscach, "Automatic adjustment of off-the-shelf reverberation effects," in *Audio Engineering Society Convention 126*, Audio Engineering Society, 2009.
- [29] A. F. Gad, *Pygad: An intuitive genetic algorithm python library*, 2021. arXiv: 2106.06158 [cs.NE].
- [30] B. Man, B. Leonard, R. King, J. D. Reiss, *et al.*, "An analysis and evaluation of audio features for multitrack music mixtures," in *Proceedings of ISMIR 2014*, 2014.
- [31] S. Hafezi and J. D. Reiss, "Autonomous multitrack equalization based on masking reduction," *Journal of the Audio Engineering Society*, vol. 63, no. 5, 2015.
- [32] D. Moffat and M. Sandler, "An automated approach to the application of reverberation," in *Audio Engineering Society Convention*, 2019.
- [33] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," 2015.
- [34] I. Smule, *DAMP-VSEP: Smule Digital Archive of Mobile Performances - Vocal Separation*, Zenodo, Oct. 2019. [Online]. Available: <https://doi.org/10.5281/zenodo.3553059>.
- [35] R. M. Bittner, J. Salamon, M. Tierney, M. Mauch, C. Cannam, and J. P. Bello, "MedleyDB: A multitrack dataset for annotation-intensive mir research," in *ISMIR*, vol. 14, 2014.
- [36] T.-S. Chan, T.-C. Yeh, Z.-C. Fan, *et al.*, "Vocal activity informed singing voice separation with the ikala dataset," in *ICASSP*, IEEE, 2015.
- [37] A. Liutkus, D. Fitzgerald, Z. Rafii, B. Pardo, and L. Daudet, "Kernel additive models for source separation," *IEEE Transactions on Signal Processing*, vol. 62, no. 16, 2014.
- [38] Harmonix, *Rock Band*, MTV Games, Electronic Arts, 2008.